

- * An artificial neuron is a simple computational model that includes input layer with weights, an activation function and an output.
- * Perceptrons
- * Input layer: It is the initial layer where input data enters the network. Each neuron in the input layer represents one feature or input variable.

Hidden layer:

Layers between the input and output layer responsible for learning complex patterns & representation from the input data. Here each neuron in a hidden layer is interconnected to neurons in the previous & subsequent layers.

Bias Value: Additional parameter to control the threshold for activation.

* Activation function: Introduce non-linearity to the model, takes weighted sum of the inputs and produces an output by capturing complex patterns and relationships.

Q Loss Function: Difference between predicted value and actual value.

(*) Cost Function: Measures the overall error all training examples.

i) for regression... MSE

ii) Binary classification... binary cross entropy

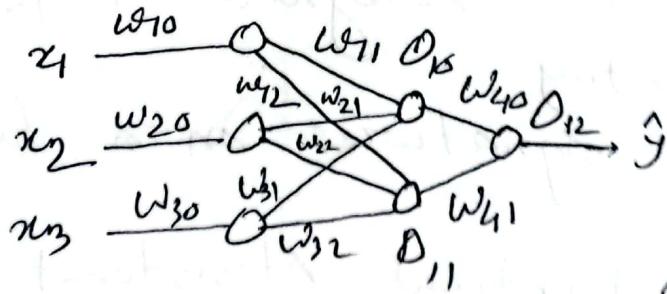
* Forward propagation:

Process of passing input data through the network layer, applying weights, biases, activation function etc. to produce predictions.

(*) Backward propagation:

Process of calculating the gradient of loss with respect to the networks parameters & use this gradient to update the parameter using optimization algorithm.

* Chain Rule: used for weight update



$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\delta L}{\delta w_{\text{old}}}$$

where : $\eta \rightarrow$ learning rate

$\frac{\delta L}{\delta w_{\text{old}}} \rightarrow$ Slope

$$\frac{\delta L}{\delta w_{412}} = \frac{\delta L}{\delta y_L}$$

weight update for: w_{411}

$$\Delta w_{411\text{new}} = w_{411\text{old}} - \eta \frac{\delta L}{\delta w_{411}}$$

$$\frac{\delta L}{\delta w_{411\text{old}}} = \frac{\delta L}{\delta o_2} \times \frac{\delta o_2}{\delta o_{11}} \times \frac{\delta o_{11}}{\delta w_{411\text{old}}}$$

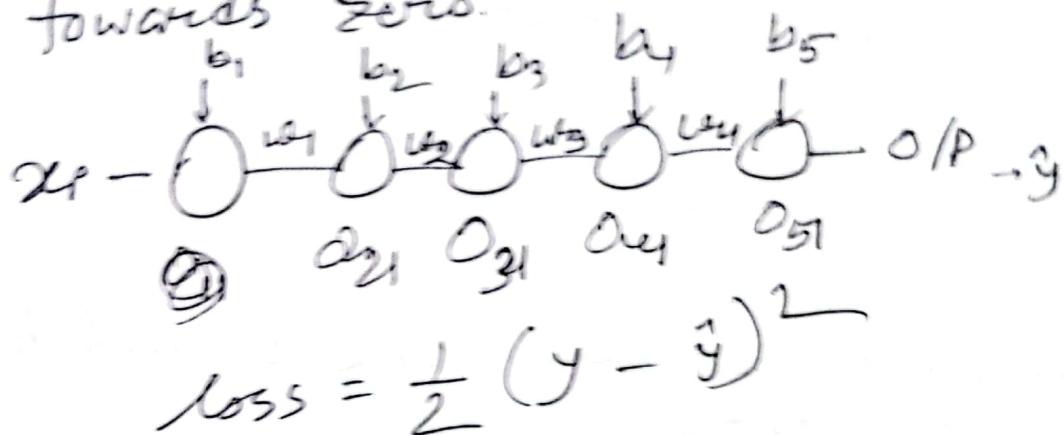
Output layer: final layer of the network
that predicted output.

Real-world application of ANN:

- i) Optical character recognition for data entry
- ii) Validation of signatures on a bank cheque.
Siri, Alexa (Virtual Assistant)
- iii)

Vanishing Gradient Descent:

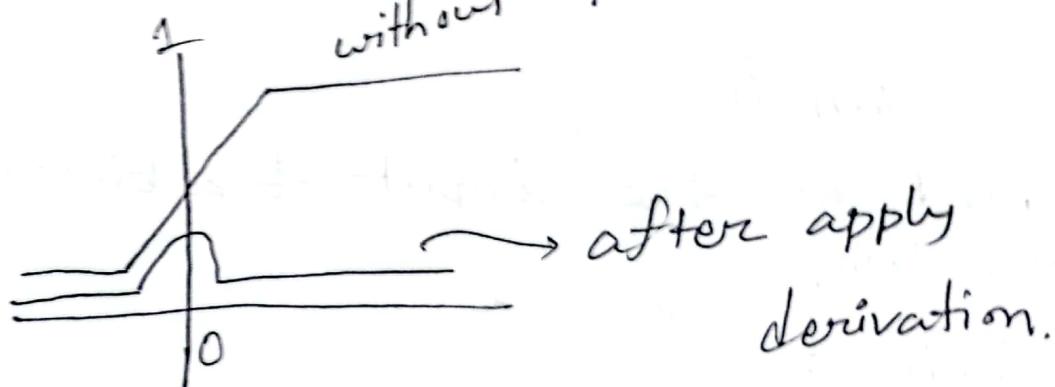
- ΔL 's becomes very small and tends towards zero.



$$O_{51} = \sigma \left[(O_{41} \times w_4) + b_5 \right]$$

↳ sigmoid

$w_{old} \approx w_{new}$



Kinds of Activation Function and their pros and cons

① Sigmoid: Normally used as the output of a binary probabilistic function.

Pros

$$f(x) = \frac{1}{1 + e^{-x}}$$

① pros

① Outputs are bounded between 0 and 1, which is useful for probabilities.

② Cons:

* Not a zero-centric function and computationally expensive.

* Tanh

* used as the input of a binary probabilistic function;

$$\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

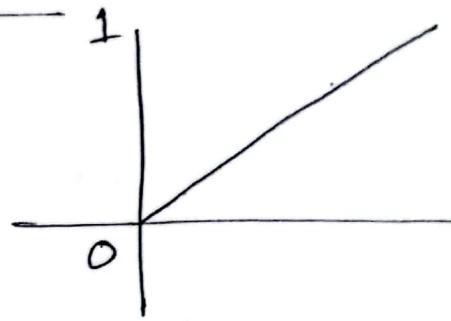
Pros

① zero centric function

cons

① prone to vanishing gradient function.

ReLU: (Rectified Linear Unit)



$$y = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

Pros

- ① Can deal with vanishing gradient problem.

Cons

- Not a zero-centric function

Softmax:

Used as the output in multi-class classification problems to find out different probabilities.

$$y = \frac{e^x}{\sum e^x}$$

Pros
Outputs sum to 1

Cons
can saturate for inputs with large absolute values.

CNN

→ CNN (Convolutional Neural Network).

Architecture:

- kernel
- stride value
- Padding
- Pooling
- Flatten

Kernel:

- Kernel is a filter, which is extracting feature from input data such as images.
- Kernel is also a matrix, which is moving onto a input and its output also a matrix dot product.

→ Formula: $\text{output} = [i - k] + i = 3$

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	2

5×5

$$\begin{matrix} 0 & 1 & 2 \\ 2 & 2 & 0 \\ 0 & 1 & 2 \end{matrix} = \begin{matrix} 3 \times 3 \end{matrix}$$

12	12	17
10	17	19
9	6	14

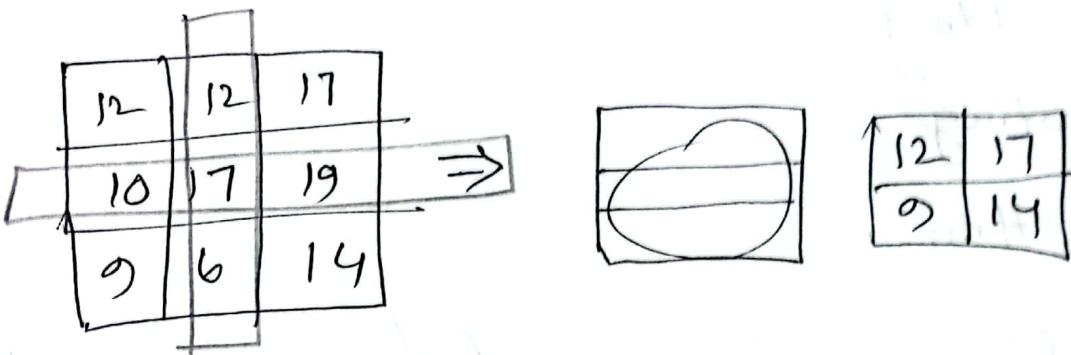
3×3

$K = \text{kernel}$

Stride value:

- Filter is moved across the image left to right, top to bottom.
- Amount of movement between application of the filter to the input image is referred to as the stride.

→



Output: $= \left[\frac{9 - K}{S} \right] + 1$

$\left[\frac{5 - 3}{2} \right] + 1 = 2$

Padding:

- Padding fixes the border effect problem.
- There is a possibility of information change due to definition change padding stores it.

1	0	4	2	3	4
8	2	5	9	3	1
6	3	4	1	1	0
0	9	3	1	0	1
1	2	3	4	5	6
7	8	9	1	2	3

6×6

$$1 \times 0 + 0 \times 1 + 4 \times 0 = 0$$

$$8 \times -1 + 2 \times 2 + 5 \times -1 = -9$$

$$6 \times 2 + 3 \times 0 + 4 \times 1 = 16$$

without padding problem:

$= 7$

① Image shrinks

② Corner pixels are less used compare the middle pixels.

padding problem

$$\begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline -1 & 2 & -1 \\ \hline 2 & 0 & 1 \\ \hline \end{array}$$

3×3

$$\begin{array}{|c|c|c|c|} \hline 7 & - & - & - \\ \hline - & - & - & - \\ \hline - & - & - & - \\ \hline - & - & - & - \\ \hline \end{array}$$

4×4

adding zero padding (same padding)

$$\text{formula} = \left\lceil \frac{i-k+2P}{s} \right\rceil + 1$$

* Pooling: Pooling are used to reduce the dimension, without losing any information.

→ It makes the model more robust to variations.

2	2	7	3
0	4	6	1
8	5	2	4
3	1	2	6

max Pool

0	7
8	6

Min Pool

2	1
1	2

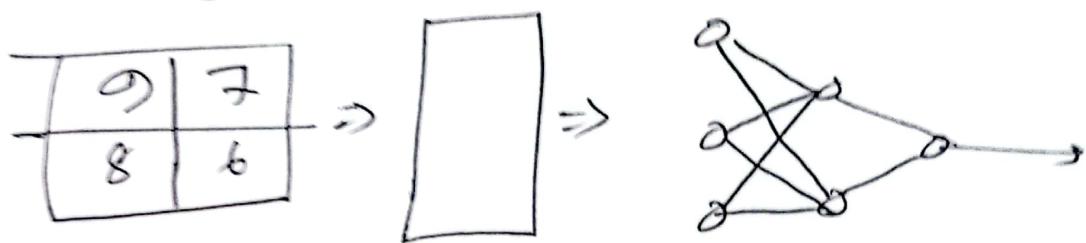
avg Pool

4.25	4.25
4.25	4.25

after
padding

Flatten:

→ Transforming the entire pooled feature map/matrix into single column which is fed to the neural network for the processing.



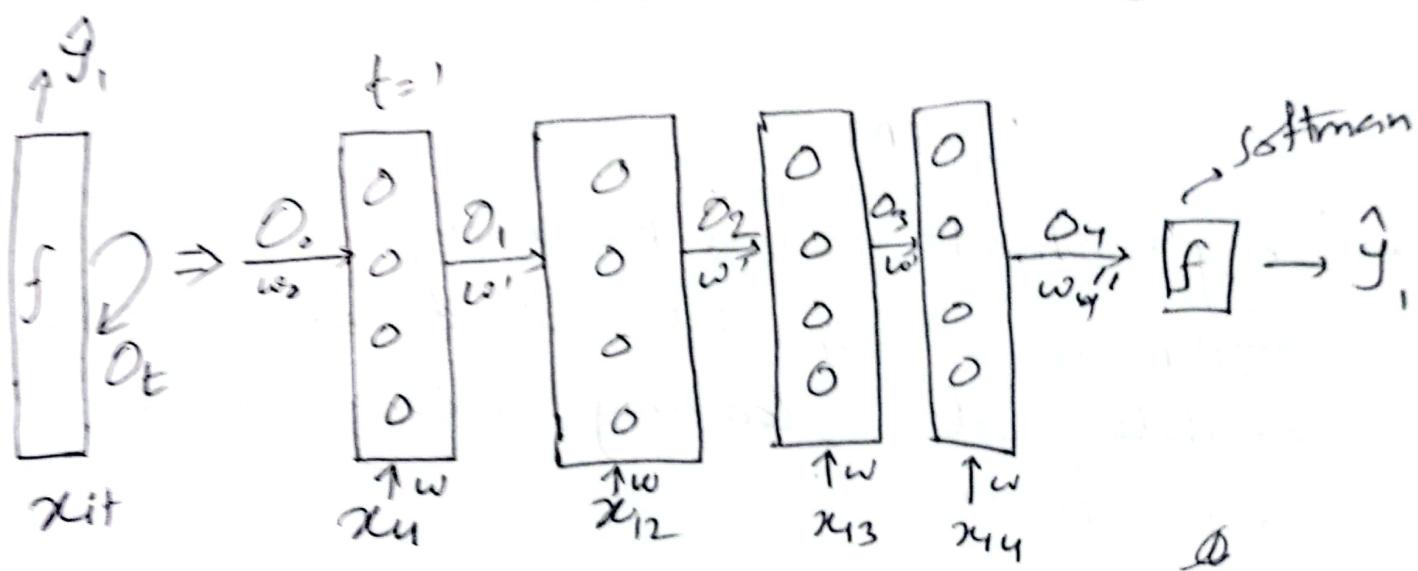
Application of CNN:

- ① Bio-metric identification of user identity
- ② Unlocking smartphones and tablets.
- ③ Airport security and passport control.

RNN

• It is specialized neural networks for handling sequential data by maintaining a hidden state that captures the context and dependencies between elements in a sequence.

↳ Application:
NLP, speech recognition.



$$O_1 = f(x_{i1} \cdot w)$$

$$O_2 = f(x_{i2} \cdot w + O_1 w')$$

$$O_3 = f(x_{i3} \cdot w + O_2 w')$$

$$O_4 = f(x_{i4} \cdot w + O_3 w')$$

Weight update:

$$w_{new} = w_{old} - \frac{\delta L}{\delta w_{old}}$$

$$\frac{\delta L}{\delta w''} = \left(\frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta w''} \right)$$

$$\frac{\delta L}{\delta w'} = \left(\frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta o_4} \times \frac{\delta o_4}{\delta w'} \right)$$

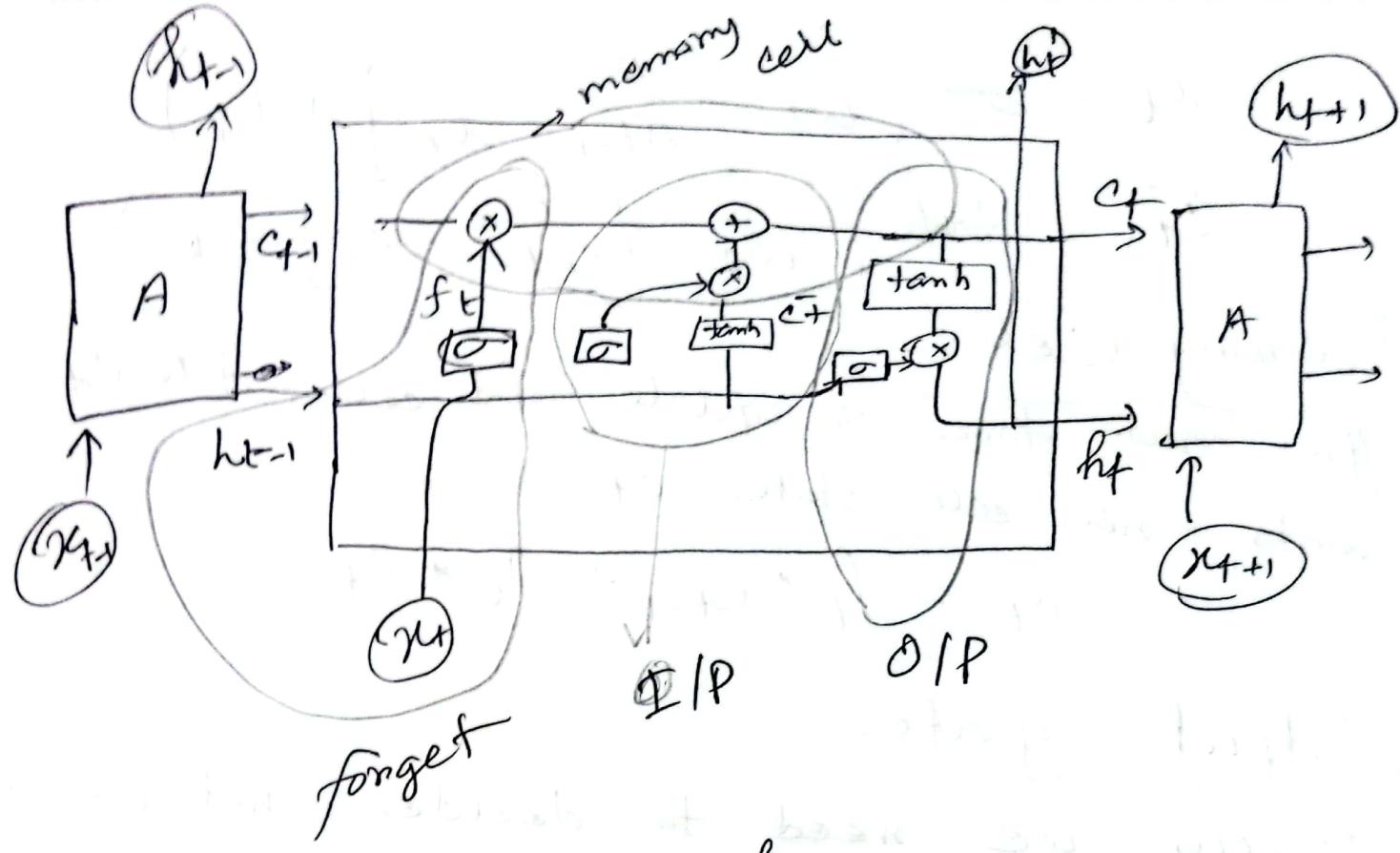
LSTM

→ Long - short term Memory

→ LSTM is a special kind of RNN capable of learning long - term dependencies.

Steps of LSTM:

- ① Memory cell
- ② Forget gate
- ③ Input gate
- ④ Output gate



① Forget gate / Sigmoid layer:

x_t and h_{t-1}

$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

when we see a new subject, we want to forget the gender of the old subject.

Input gate:

We'd want to add the gender of the new subject to the cell state, to replace the old one we were forgetting.

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

$$c_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

Memory Cell:

It's now time to update old cell state c_{t-1} ,
into new cell state c_t

$$c_t = f_t * c_{t-1} + i_t * \bar{c}_t$$

Output gate:

Finally, we need to decide what we are
going to output. First sigmoid layer
and put the cell state through the tanh
and multiply it by the output of the
sigmoid gate.

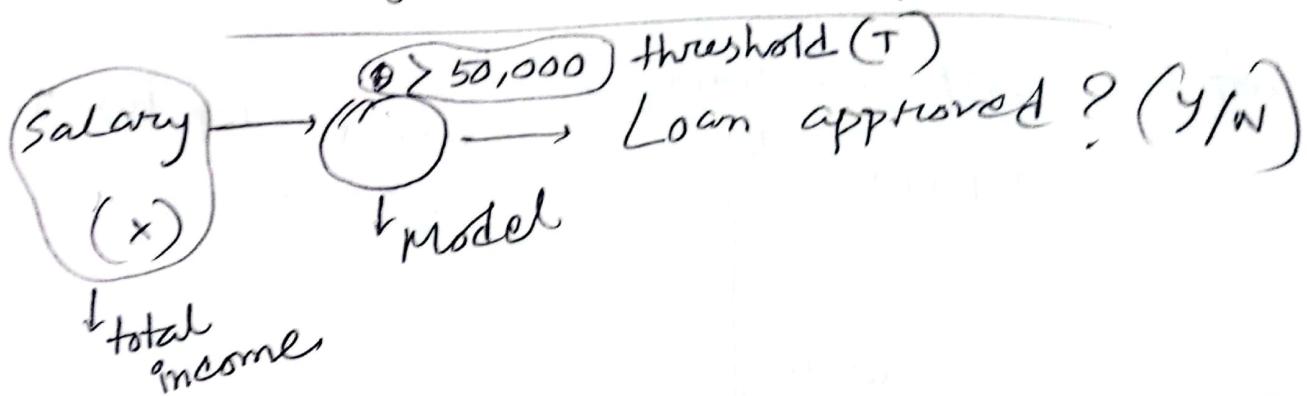
$$o_t = \sigma(w_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(c_t)$$

LSTM applications

- ① Speech recognition
- ② Rhythm learning
- ③ Music composition
- ④ Grammar learning.

Understand Perceptron



Total income = Applicant Salary (x_1) + Father salary (x_2)
+ family member (x_3)

$$\text{Total Income} = x_1 + x_2 + x_3 > 0 \quad T$$

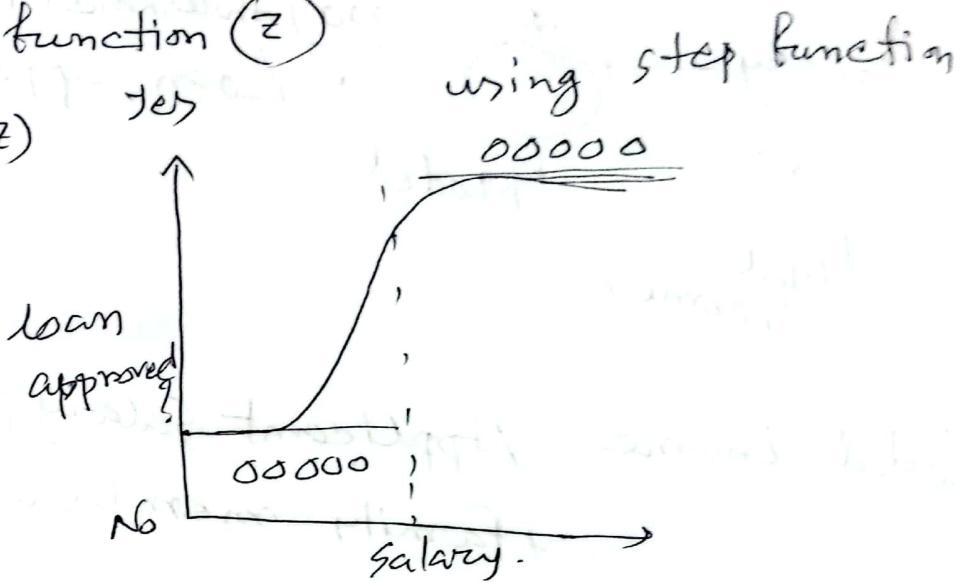
$$\Rightarrow x_1 + x_2 + x_3 - T > 0$$

$$\Rightarrow x_1 + x_2 + x_3 + \text{bias} > 0$$

If $x_1 + x_2 + x_3 + b_i > 0 \rightarrow \text{Output } 1$
 $x_1 + x_2 + x_3 + b_i \leq 0 \rightarrow \text{Output } 0$

Use of step-up function:

$$\begin{aligned} \text{Output} &= \text{step function}(z) \\ &= \text{Step}(z) \quad \text{Yes} \end{aligned}$$



Perceptron Learning Algorithm:

→ If correct: Do nothing if the prediction input is equal to the target output.

→ If incorrect:

↳ scenario: (a) if output is 0 and target 1, add input vector to weight vector.

↳ scenario: (b) if output is 1 and target 0, subtract input vector from weight vector.

let,

$$D = \left((x^{[1]}, y^{[1]}), (x^{[2]}, y^{[2]}), \dots, (x^{[n]}, y^{[n]}) \right) \in \left(\mathbb{R}^m \times \{0, 1\} \right)^n$$

A. for every $(x^{[i]}, y^{[i]}) \in D$:

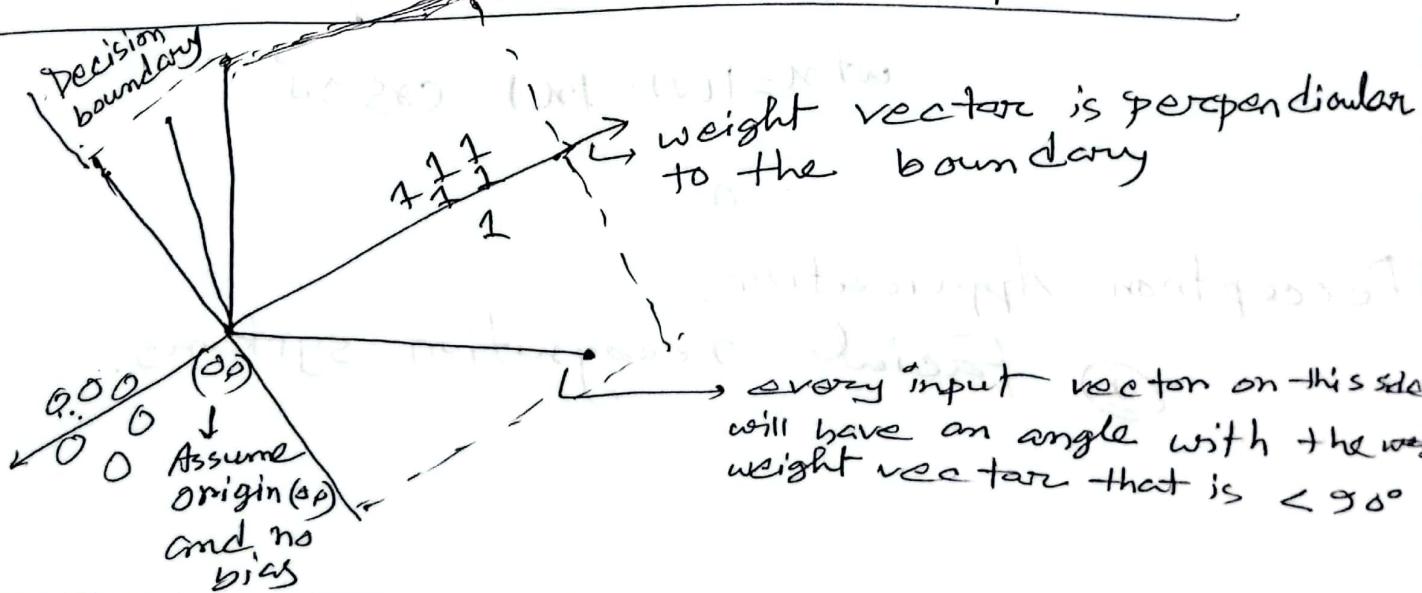
① $\hat{y}^{[i]} := \sigma(x^{[i] \top} w)$

② error := $(y^{[i]} - \hat{y}^{[i]})$

③ $w_{\text{new}} := w_{\text{old}} + \text{err} \times x^{[i]}$

↳ weight update.

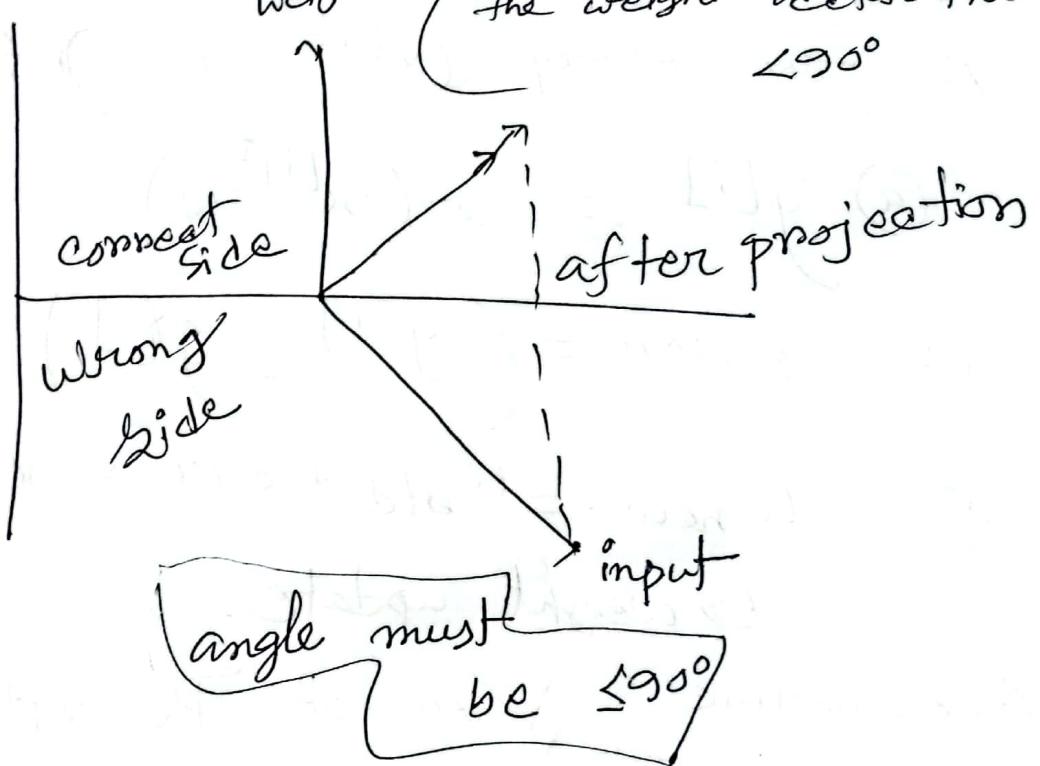
Geometric View of Perceptron:



$$\hat{y} = \begin{cases} 0 & w^T x \leq 0 \\ 1 & w^T x > 0 \end{cases}$$

$\cos 90^\circ = 0$

$w^T x = \|w\| \cdot \|x\| \cdot \cos \theta$ so this needs to be 0 at the boundary if it will have an angle with the weight vector that is 90°



$$w^T x = \|w\| \|x\| \cos 90^\circ$$

$$\leq 0$$

Perception Applications:

- ① Facial recognition systems.